

**Achtung!**

Dies ist eine Internet-Sonderausgabe des Aufsatzes "Perspektiven des Computereinsatzes in der Orientalistik" von Jost Gippert (1989).

Sie sollte nicht zitiert werden. Zitate sind der Originalausgabe in "Forschungsforum" 2, 1990, 133-136 zu entnehmen.

**Attention!**

This is a special internet edition of the article "Perspektiven des Computereinsatzes in der Orientalistik" ["Perspectives of the application of computers in Near Eastern studies"] by Jost Gippert (1990).

It should not be quoted as such. For quotations, please refer to the original edition in "Forschungsforum" 2, 1990, 133-136.

**Alle Rechte vorbehalten / All rights reserved:**

Jost Gippert, Frankfurt 1999

# Perspektiven des Computereinsatzes in der Orientalistik

Jost Gippert

Nachdem die Domäne des Computers im universitären Einsatz noch bis vor wenigen Jahren im naturwissenschaftlichen Bereich lag, werden elektronische Verfahren heute mehr und mehr auch in den Geisteswissenschaften angewendet. Das primäre Einsatzgebiet liegt dabei zweifellos in der Textverarbeitung; der Computer erweist sich hier als ein universal einsetzbares Hilfsmittel, das die Gestaltung eines Textes von seiner Konzipierung bis zur Drucklegung in der Hand des Autors ermöglicht und herkömmlichen Verfahren somit überlegen ist.

Als entscheidende Voraussetzung für einen weitergehenden Siegeszug des Computers in den Geisteswissenschaften ist die Verfügbarkeit von Zeichen und Schriften anzusehen, die über den Grundvorrat des lateinischen Alphabets hinausgehen. Während im Bereich von Großrechenanlagen und im internationalen Datentransfer noch immer ein Zeichenvorrat von nur 128 Zeichen verwendet wird, der in der sog. ASCII-Norm die 26 Buchstaben des Lateinalphabets in Groß- und Kleinschreibung, die Ziffern von 0 bis 9, die wichtigsten Satzzeichen sowie 32 Steuercodes für Zeilenschaltungen u.ä. umfaßt, war schon die Aufstockung auf 256 Zeichen, mit denen heute jeder Personal Computer nach dem IBM-System operiert, ein bedeutender Fortschritt. Dieser erweiterte Zeichensatz ist auf die west- und nordeuropäischen Nationalalphabete deutsch, französisch, italienisch, spanisch, niederländisch und schwedisch zugeschnitten und umfaßt über die 128 ASCII-Zeichen hinaus z.B. diakritische Kombinationen wie *ä, ö, ü, á, é, â* etc., die Ligatur *æ*, die Währungssymbole für das engl. Pfund (£) oder den holl. Gulden (f), einige in der Mathematik gebrauchte Symbole wie z.B.  $\cap$  oder  $\equiv$ , einige - ebenfalls nach den Bedürfnissen der Mathematik ausgewählte - griechische Buchstaben wie z.B.  $\alpha$  oder  $\Sigma$  sowie sog. Grafikzeichen wie z.B.  $\lfloor$  oder  $\lceil$ , die zur Erzeugung von Linien oder Rahmen am Bildschirm gebraucht werden. Als ein Kuriosum ist festzuhalten, daß das deutsche "scharfe" *ß* in diesem Zeichensatz fehlt; stattdessen muß das Zeichen *β* verwendet werden, das im System das griech. *beta* repräsentiert.

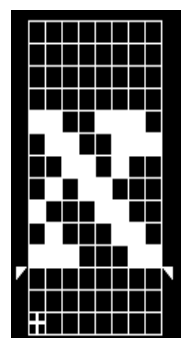
Es liegt auf der Hand, daß für einen universalen Einsatz in den Geisteswissenschaften auch dieser erweiterte Zeichensatz noch keine geeignete Grundlage abgibt. Weder ein Slavist, der z.B. tschechisches Wortmaterial zu verarbeiten hat, noch ein Klassischer Philologe, der das Altgriechische nicht transkribieren will, findet sämtliche von ihm benötigten Zeichen wieder; von den Bedürfnissen eines Orientalisten, der mit Schriften wie Hebräisch oder Arabisch arbeiten muß, ist dabei ganz zu schweigen. Diese Probleme sind innerhalb der letzten fünf Jahre an verschiedenen Orten und mit unterschiedlichen Verfahren angegangen worden, und

zufriedenstellende Lösungen liegen bereits in großem Umfang vor. Das betrifft zunächst die Erweiterung des Zeichenvorrats um weitere lateinschriftliche Sonderzeichen und diakritische Kombinationen, dann die Verfügbarkeit von kompletten, auf das Altgriechische zugeschnittenen griechischen Zeichensätzen einschließlich der zahlreichen Akzentkombinationen sowie das kyrillische Alphabet, wie es heute im Russischen und den anderen slavischen Nationalsprachen verwendet wird. Im Bereich der Orientalistik bleibt jedoch noch genügend Entwicklungsarbeit zu leisten, bis die hier anfallenden Sprachen und Schriften dem Benutzer eines Computers mit demselben Komfort zur Verfügung stehen. Was bei einer solchen Entwicklungsarbeit zu bedenken ist, soll im folgenden kurz umrissen werden.

Die Anforderungen, die an ein mit Originalschriften operierendes Textverarbeitungssystem gestellt werden, müssen sich prinzipiell an dem heute erreichten Leistungsstandard der "normalen" lateinschriftlichen Textverarbeitung orientieren. Das bedeutet zunächst, daß die betreffende Schrift am Bildschirm sichtbar sein muß und, möglichst in verschiedenen Schriftgrößen und -ausgestaltungen, auch auf dem Drucker ausgegeben werden kann, ferner, daß sie in einer adäquaten Weise, im Normalfall über die Tastatur, erzeugbar sein muß. Es bedeutet weiter, daß die Standardfunktionen der Textverarbeitung wie z.B. automatisches Suchen und Ersetzen bestimmter Zeichen oder Zeichenfolgen, Fuß- und Endnotenverwaltung, Zeilenfunktionen wie Blocksatz, Zentrieren oder rechtsbündige Ausrichtung, Spalten- und Tabellenerstellung in der Fremdschrift mit derselben Effektivität ausgeführt werden können wie in der Lateinschrift. Wünschenswert ist darüber hinaus eine eindeutige Kodierung, die es ermöglicht, den erzeugten Text von einem System auf ein anderes zu übertragen; dies ist z.B. die Voraussetzung für die Weitergabe an eine Druckerei ("Drucken von Diskette"). Eine unabdingbare Anforderung, die die wissenschaftliche Arbeit stellt, ist letztlich die gleichzeitige Verfügbarkeit von Originalschrift(en) und Lateinschrift innerhalb desselben Texts.

Um diesen Erfordernissen gerecht zu werden, bedarf es in der Regel eines vielschichtigen Systems von ineinandergreifenden Programmelementen, sog. Treibereinheiten, die von der internen Struktur des verwendeten Rechners und seiner Zusatzgeräte abhängen. Das betrifft zum einen die Erzeugung von Schriftzeichen auf Bildschirm und Drucker, die vorrangig auf sogenannten Bitmatrizen, d.h. Rastern aus Einzelpunkten, basiert. Umfang und Dichte dieser Raster hängen einerseits von der Leistungsfähigkeit der verwendeten Geräte ab, andererseits von der gewünschten Größe der Zeichen. Eine beispielhafte Zusammenstellung ist den folgenden Abbildungen zu entnehmen, wo der hebräische Buchstabe Aleph in fünf Matrizes dargestellt ist: Abbildung 1 zeigt das Aleph in einer typischen Matrix für die Bildschirmausgabe,

Abb. 1: Bitmatrix für Bildschirmzeichensatz



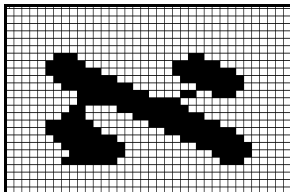


Abb. 2: Bitmatrix für 24-Nadel-Drucker (24 · 36 Punkte)

die nach dem sogenannten "EGA" - Standard ("Enhanced Graphics Adapter") 14 · 8 Punkte umfaßt. Abbildung 2 zeigt das Aleph in der Bitmatrix für 24-Nadel-Drucker, die im Schönschriftmodus üblicherweise mit einer Matrix von 24 · 36 Punkten operieren. Abbildungen 3 bis 5 zeigen denselben Buchstaben, wie er in den drei Größen 8-Point, 10-Point und 12-Point (ein "Point" = 1/72 Zoll) auf einem handelsüblichen Laserdrucker mit einer Punktdichte von 300 · 300 Punkten pro Quadratzoll ausgegeben werden kann; im Original entsprechen diese drei Matrizes etwa den Druckgrößen  $\aleph$ ,  $\aleph$  und  $\aleph$ .

Die Erstellung der notwendigen Zeichensätze (sog. "Fonts") ist also der erste zu bewältigende Schritt bei der Entwicklung. Allerdings können die verschiedenen Größen nur unter großem Qualitätsverlust automatisch auseinander abgeleitet werden. Hier zeichnet sich für die Zukunft eine Vereinfachung ab, insofern für Laserdrucker anstelle von Bitmatrix-

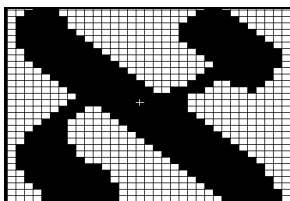
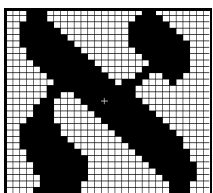
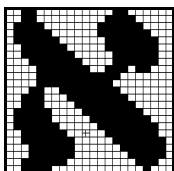


Abb. 3 bis 5: Bitmatrizes für Laserdrucker (8-Pt. / 10-Pt. / 12-Pt.-Größe)

fonts mehr und mehr auf sogenannte Outlinefonts übergegangen wird, bei denen nicht mehr einzelne Punkte, sondern die Umrisse der Zeichen als Linien mit Anfangs- und Endpunkten definiert werden, die der Drucker durch Umrechnung dann selbständig auf die gewünschte Größe bringt (vgl. Abbildung 6 mit einer Outline-Darstellung des hebräischen Aleph). Sowohl für die Generierung über Matrizes als auch für die von Outlinefonts stehen heute kommerzielle Programme zur Verfügung.

Eine zweite Entwicklungseinheit betrifft die Verwaltung der erstellten Zeichen durch den Rechner. Hier gibt es grundsätzlich zwei konkurrierende Verfahren, den sog. Text- oder Alphamodus und den sog. Grafikmodus. Beide unterscheiden sich im wesentlichen dadurch, daß im ersteren Fall zunächst der vollständige Zeichensatz in den für den Bildschirm vorgesehenen Speicher (sog. "Bildschirm- oder Grafikkarte") bzw. in den Speicher des Druckers geladen wird (sog. "Downloadverfahren") und die Geräte daraufhin über bestimmte Codes angewiesen werden, die Matrix des anzusteuern Zeichen aus diesem Zeichensatz herauszulesen und auszugeben; im Grafikmodus wird hingegen die erforderliche Punktmatrix für jedes einzelne darzu-

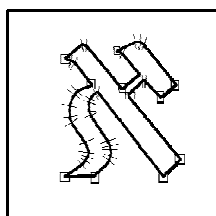
stellende Zeichen vom Rechner an die Ausgabegeräte weitergeleitet. Das letztere Verfahren ist wesentlich flexibler, da die Anzahl der gleichzeitig verwendbaren Zeichensätze hier praktisch unbegrenzt ist, während sie im Downloadverfahren von der Kapazität des Druckers bzw. der Bildschirnkarte abhängt. Andererseits erlaubt der Textmodus im Normalfall eine wesentlich schnellere Verarbeitung, da der für die Definition eines Zeichens benötigte, bei größeren Matrizes recht umfangreiche Code nur einmal ausgegeben zu werden braucht und dann über einen wesentlich kürzeren Code abrufbar ist. Die Verwaltung der Zeichensätze und der Steuercodes ist die zentrale Funktion der sog. Bildschirm- und Druckertreiber.

Eine am gestalterischen Optimum orientierte Verarbeitung von Schriften wird nicht nur bei Lateinschriften, sondern auch bei den meisten anderen Schriften von den bei Schreibmaschinen üblichen festen Schrittabständen abgehen und stattdessen eine "proportionale" Wiedergabe vorziehen, bei der jedem einzelnen Zeichen die seiner Form entsprechende Breite zukommt; man vgl. etwa die Lateinbuchstaben *i* und *m*, die sich in ihrer natürlichen Breite etwa um einen Faktor drei unterscheiden. Einem Programm, das die optimale Füllung einer Zeile oder auch die nötigen Wortzwischenräume für die Erstellung von Blocksatz errechnen soll, muß die spezifische Weite eines jeden Zeichens bekannt sein; diese Aufgabe übernehmen die sog. "Weitentabellen", die eine eigene zu erstellende Treibereinheit bilden.

Die wohl aufwendigste Treibereinheit betrifft die Ansteuerung der gewünschten Zeichen über die Tastatur. Prinzipiell ist davon auszugehen, daß jede Taste auf der Tastatur eines Computers beim Niederdrücken einen bestimmten elektronischen Code erzeugt, der vom Computer empfangen und interpretiert wird. Daß beim Niederdrücken der Taste "M" tatsächlich ein "M" erzeugt wird, ist nicht selbstverständlich und beruht ausschließlich auf der Interpretation des Tastencodes durch den Rechner. Durch ein eigenes Programm kann der Rechner nun angewiesen werden, aufgrund desselben Tastencodes einmal ein "M" und einmal ein arabisches  $\aleph$  zu erzeugen. Je nach der Menge und Art der zu verarbeitenden Einzelzeichen einer Schrift und ihrer internen Codierung durch das Textverarbeitungsprogramm ist die Gestaltung eines solchen Tastaturtreibers mehr oder weniger komplex, wobei im orientalischen Bereich die Sonderproblematik von linksläufigen Schriften wie der arabischen, vertikal angeordneten Schriften wie der klassisch-mongolischen, von Silbenschriften wie der koreanischen oder von Wortschriften wie der chinesischen zu berücksichtigen ist. Pauschale Lösungen gibt es hier nicht; vielmehr benötigt jede Schrift ihre eigene Anpassung, die vom Ideal eines möglichst einfachen Eingabe- und Umschaltmodus geprägt sein sollte.

Anzustreben ist nach alledem eine Textverarbeitung, die den gesamten Bereich orientalischer Schriften abdeckt. Die Realisierung eines solchen Vorhabens, an der die Universität Bamberg durch die neugeschaffene Arbeitsstelle für Orientalistische Computerlinguistik beteiligt ist, wird noch einige Zeit in Anspruch nehmen.

Abb. 6: Hebräisches Aleph in Outline-Darstellung





Die meisten anderen orientalischen Schriften wie z.B. die arabische verlangen demgegenüber ein abweichendes Verfahren, bei dem eine andere Abgrenzungsbasis als der Abstand zwischen einzelnen Zeichen zu wählen ist; die Entwicklung solcher Verfahren ist als ein vordringliches Desiderat einzustufen.

Ein weiteres Desiderat im Hinblick auf den Einsatz des Computers bei der sprachwissenschaftlichen Analyse orientalischer Textmaterialien sind Verfahren für eine morphologische Sortierung. Obwohl schon heute durchaus leistungsfähige Programme vorliegen, die Wortformenindizes und Textstellenkonkordanzen aus elektronisch gespeicherten Texten ableiten, reichen die erzielten Resultate für wissenschaftliche Fragestellungen häufig nicht aus, da das Ordnungsprinzip im Normalfall eine Auflistung nach dem Alphabet der Wortformen ist; um ein benutzbares Lexikon zu schaffen, wird hingegen eine Lemmatisierungsfunktion benötigt, die von suffixalen oder sogar präfialen Elementen abstrahiert und Flexionsformen in ihren paradigmatischen Zusammenhang einordnet. Man vergleiche dazu z.B. die folgenden Auszüge aus einem automatisch erstellten Index und einer ebenso erzeugten Konkordanz zu dem etwa 300 Seiten umfassenden Sammelband von Volksliedtexten in der Sprache der kaukasischen Svanen (Svanuri Poezia, Tbilisi 1939):

Taf. 4: Wortformenindex zu Svanuri Poezia (Auszug)

Wortformenindex zu Svanuri Poezia:		
(In Klammern: Die Häufigkeit der betreffenden Wortform)		
<i>abaz</i>	(1)	110h: 350, 4
<i>abram</i>	(2)	64b: 236, 82; 66: 242, 25
<i>abrams</i>	(1)	66: 240, 3
<i>abrešumiš</i>	(1)	1b: 4, 13
<i>abrešvimiš</i>	(1)	50a: 174, 24
<i>abžinalix</i>	(1)	5: 18, 46
<i>abžare</i>	(1)	1a: 2, 7
<i>abžari</i>	(1)	29: 100, 7
<i>abžaris</i>	(2)	32: 110, 14; 41b: 136, 14
<i>abžriš</i>	(1)	41b: 138, 52
<i>acar</i>	(1)	66: 242, 12
<i>acars</i>	(2)	66: 240, 7; 242, 15
<i>aceri</i>	(1)	66: 242, 20
<i>acurax</i>	(1)	94a: 290, 36
<i>acvir</i>	(2)	31: 108, 65; 43b: 154, 64
<i>acvird</i>	(1)	20: 66, 37
<i>ačanyeli</i>	(1)	30: 102, 2
<i>ači</i>	(1)	67: 244, 27
<i>ačad</i>	(18)	8: 30, 102; 18: 62, 13; 27a: 92, 73; 41a: 132, 9; 41b: 136, 9.11; 43b: 152, 39; 77a: 256, 1.3.4.5.6. 7.8; 258, 11.12; 77b: 258, 4; 91b: 270, 18
<i>ačadd</i>	(1)	57a: 190, 37
<i>ačadx</i>	(2)	9: 38, 56; 94a: 294, 118
<i>ače</i>	(2)	28: 98, 22.33
<i>ačed</i>	(1)	39b: 124, 5
<i>ačkad</i>	(2)	22: 70, 15.17
<i>ačungo</i>	(1)	26: 88, 72

#### Textstellenkonkordanz zu Svanuri Poezia:

		<i>xoqaci</i> (1)
72/Mlx:	(248),6	<i>doşgu xoqaci macxvari!</i>
		<i>xoqde</i> (1)
43a/Lšx:	(150),30	<i>merma katxas gvalvars xoqde".</i>
		<i>xoqdex</i> (2)
25b/Mlx:	(80),25	<i>xexvas mineš ečav xoqdex;</i>
	(80),27	<i>begärs mine šonšv ečav xoqdex;</i>
		<i>xoqida</i> (4)
8/Mlx:	(28),79	<i>päsild utkläbvd ka xoqida.</i>
9/Mst:	(38),50	<i>žveg i mežveg mäg xoqida.</i>
41b/Ƙal:	(134),7	<i>baṗəld ägite xoqida,</i>
46/Mlx:	(166),107	<i>mišgov lerekv ka xoqida;</i>
		<i>xoqidax</i> (5)
8/Mlx:	(26),54	<i>atxa känte čur xoqidax,</i>
13/Mlx:	(46),15	<i>sga xoqidax, qän sga xobax,</i>
23/Lžr:	(70),5	<i>ži xoqidax zagärteži.</i>
51/Ltl:	(164a),68	<i>sgav xoqidax Qaräštësga,</i>
72/Mlx:	(248),2	<i>žvegi xoqidax didi mindvriše.</i>

Taf. 5: Textstellenkonkordanz zu Svanuri Poezia (Auszug)

Obwohl auch solche rein alphabetischen Auflistungen bereits einen eigenen Wert haben, da sie als Grundstock für eine vollständige Erfassung der auftretenden Wortformen und darauf aufbauende morphologische und phonologische Analysen dienen können, wäre eine weitergehende Nutzung erst dann möglich, wenn die paradigmatische Zusammengehörigkeit von Wortformen wie *ačad* und *otčed* (beides finite Formen einer Verbalwurzel *-čed-* "gehen") oder von *xoqde* und *miqida* (beides finite Formen einer Verbalwurzel *-qed-* "gehen") erkannt würde und als übergelagertes Ordnungsprinzip eingesetzt wäre. Eine solche Funktion muß natürlich die internen grammatischen Regeln der jeweiligen Objektsprache reflektieren und kann nicht ohne weiteres auf pauschale Algorithmen zurückgreifen. Auch hierzu bedarf es weiterführender Entwicklungen, an denen sich die Arbeitsstelle Orientalistische Computerlinguistik der Universität Bamberg in Kooperation mit anderen in- und ausländischen Hochschulen beteiligen wird.